

# Support Vector Machines

## Finding the Optimal Decision Boundary

Ahmed BADI

Mathematics & Machine Learning Enthusiast

*ahmedbadi905@gmail.com*  
*linkedin.com/in/badi-ahmed*

December 15, 2025

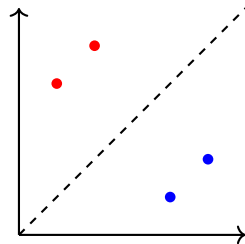
# Overview

- 1 Introduction
- 2 The Geometric Perspective
- 3 Mathematical Foundation
- 4 The Dual Formulation
- 5 Soft Margin SVM
- 6 The Kernel Trick
- 7 Extensions & Practical Considerations

# Introduction

# The Intuition: Apples and Oranges

- **The Goal:** Separate two classes (e.g., apples vs. oranges).
- **The Question:** Which line is the best?
- **Intuition:** We want the "widest road" possible between the classes.
- Maximizing the distance to the nearest points gives us more "breathing room" for new data.



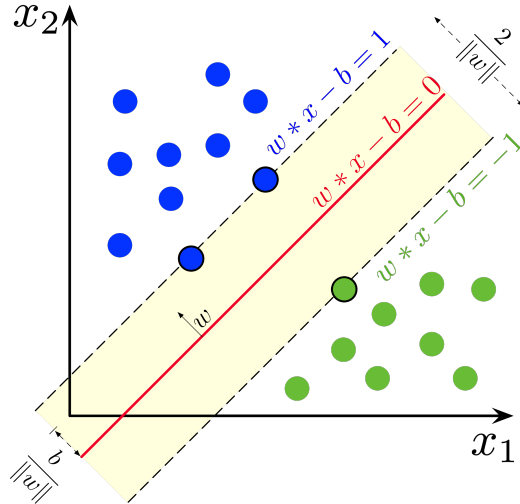
# What makes SVMs special?

## Key Characteristics

- **Optimality:** Finds the *best* boundary, not just *any* boundary.
- **Theoretical Foundation:** Rooted in Statistical Learning Theory (Vapnik-Chervonenkis dimension).
- **Sparsity:** The solution depends only on a few key data points (**Support Vectors**).
- **Versatility:** Can handle non-linear patterns via the **Kernel Trick**.

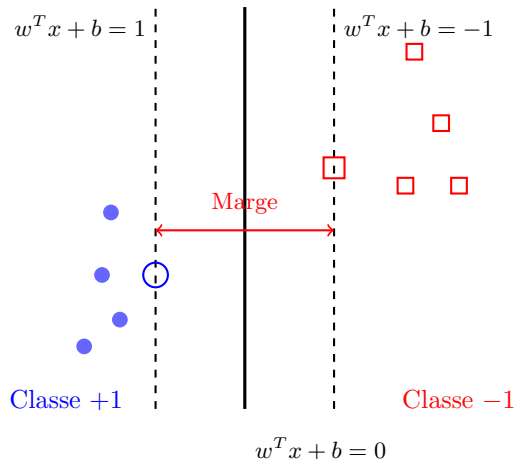
# The Geometric Perspective

# Which hyperplane is best?



# The Margin Concept

- **Margin:** Distance from decision boundary to nearest data point.
- **Support Vectors:** The critical points lying on the margin boundaries.
- Larger margin  $\Rightarrow$  Better generalization & robustness.





# Mathematical Foundation

# Problem Setup

- **Data:**  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, +1\}$ .
- **Hyperplane:**  $w^T x + b = 0$ .
- **Decision Function:**  $f(x) = \text{sign}(w^T x + b)$ .

## Geometric Margin

The distance between the marginal hyperplanes is:

$$\text{Margin} = \frac{2}{\|w\|}$$

# The Optimization Problem (Primal)

To maximize the margin  $\frac{2}{\|w\|}$ , we minimize  $\|w\|$ .

## Hard Margin SVM Optimization

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{subject to} \quad & y_i(w^T x_i + b) \geq 1, \quad \forall i = 1, \dots, n \end{aligned}$$

- Convex quadratic programming problem.
- Unique global minimum exists.

# The Dual Formulation

# Why the Dual?

Using Lagrange multipliers  $\alpha_i$ , we transform the problem.

- **Weight Vector:**  $w = \sum_{i=1}^n \alpha_i y_i x_i$
- $w$  is a linear combination of training points.

## Key Property

- If  $\alpha_i = 0$ : Point is correctly classified, not on margin (irrelevant).
- If  $\alpha_i > 0$ : Point is a **Support Vector**.

# The Dual Optimization Problem

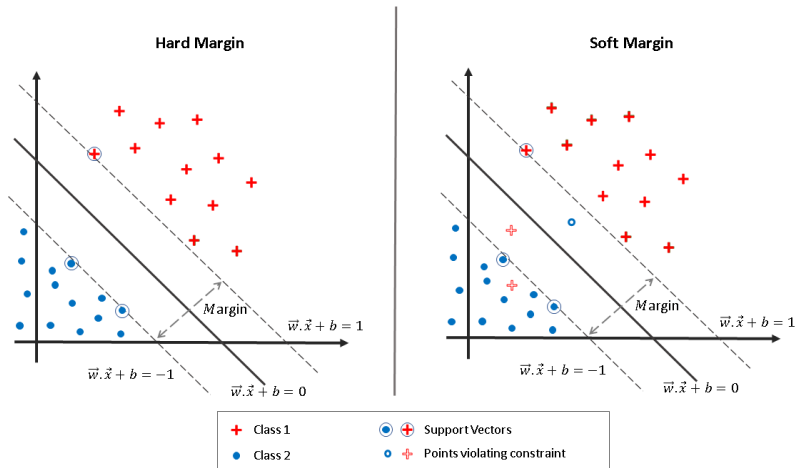
$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{subject to} \quad & \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

**Crucial Observation:** The problem depends *only* on dot products  $\mathbf{x}_i^T \mathbf{x}_j$ . This enables the Kernel Trick.

# Soft Margin SVM

# Handling Non-Linearly Separable Data

Real data has noise and outliers. A "Hard Margin" is often impossible.





# Slack Variables and Parameter C

We introduce slack variables  $\xi_i \geq 0$  to allow margin violations.

## Optimization with C

Minimize:  $\frac{1}{2} \|w\|^2 + C \sum \xi_i$

### Role of C:

- **Large C:** Penalizes errors heavily (low bias, high variance).
- **Small C:** Wider margin, tolerates errors (smoother boundary).

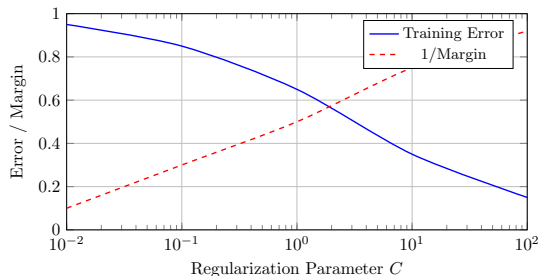
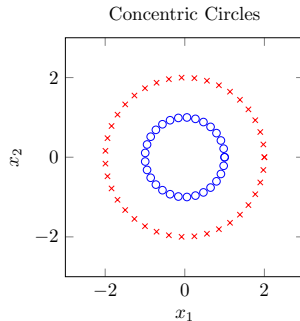
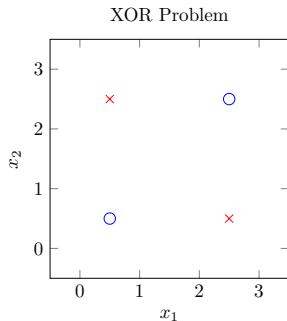


Figure: Effect of C on the boundary.

# The Kernel Trick

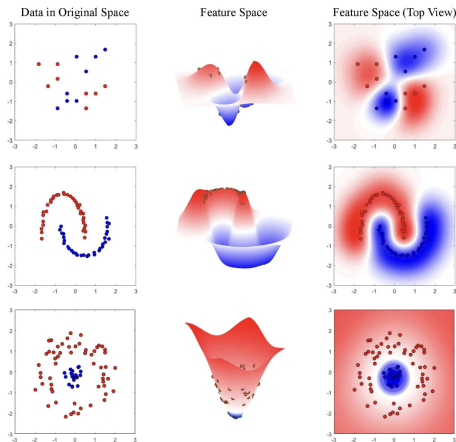
# Limitations of Linear Boundaries

Some problems (like XOR or concentric circles) are impossible for linear classifiers.



# Feature Mapping

- **Idea:** Map data  $x$  to a higher-dimensional space  $\phi(x)$  where it becomes linearly separable.
- **Problem:** Computing  $\phi(x)$  explicitly is expensive.



# The Trick

We replace the dot product  $\mathbf{x}_i^T \mathbf{x}_j$  with a Kernel function:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

We compute the "similarity" in high-dimensional space without ever transforming the data explicitly.

## Common Kernels:

- **Linear:**  $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$
- **Polynomial:**  $K(\mathbf{x}, \mathbf{z}) = (\gamma \mathbf{x}^T \mathbf{z} + r)^d$
- **RBF (Gaussian):**  $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$

# Kernel Examples

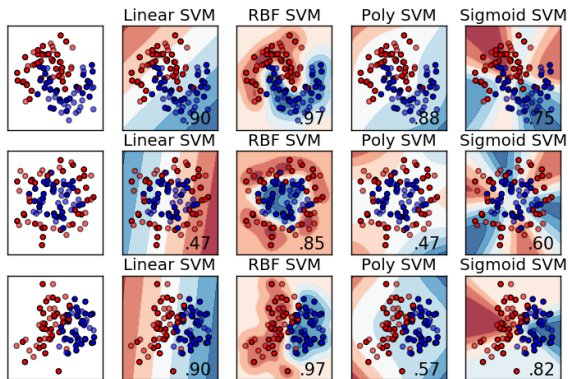


Figure: Different kernels create different decision boundaries.

## Extensions & Practical Considerations

# Multi-Class Classification

SVMs are naturally binary. For  $K$  classes:

## One-vs-Rest (OvR)

- Train  $K$  classifiers.
- Class  $k$  vs. All others.
- Faster, good for large datasets.

## One-vs-One (OvO)

- Train  $K(K - 1)/2$  classifiers.
- Pairwise competition.
- Voting system for prediction.



# Support Vector Regression (SVR)

- Fits a "tube" of width  $\epsilon$  around data.
- Points inside the tube have zero loss.
- Points outside become support vectors.

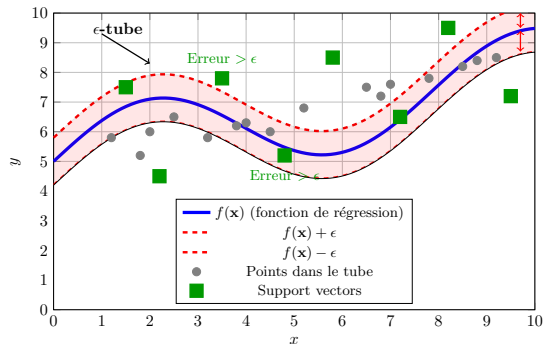


Figure: SVR with  $\epsilon$ -tube.

# Practical Tips

## Feature Scaling is Mandatory!

SVMs rely on distances. Unscaled features ruin performance.

$$x_{new} = \frac{x - \mu}{\sigma}$$

## Hyperparameter Tuning:

- **C**: Regularization.
- **Kernel**: RBF is a good default.
- **$\gamma$  (Gamma)**: For RBF, controls influence radius.

## Conclusion

# Advantages vs. Limitations

## Advantages

- High accuracy in high dimensions ( $d > n$ ).
- Memory efficient (only stores SVs).
- Global optimum (Convex).
- Versatile (Kernels).

## Limitations

- Slow training on large datasets  $O(n^2)$ .
- Sensitive to noise/overlapping.
- Requires scaling.
- Harder to interpret than Trees.

# Summary

*"The best model is not always the most complex one, but the one that generalizes best."*

SVMs provide a rigorous, elegant framework for finding optimal boundaries, remaining a fundamental tool in the Machine Learning arsenal.

# Thank You!

Questions?